

Robust Large-Scale Non-Negative Matrix Factorization Using Proximal Point Algorithm

Jason Gejie Liu and Shuchin Aeron

Department of Electrical and Computer Engineering

Tufts University, Medford, MA 02155

Gejie.Liu@tufts.edu, shuchin@ece.tufts.edu

Abstract

A robust algorithm for non-negative matrix factorization (NMF) is presented in this paper with the purpose of dealing with large-scale data, where the separability assumption is satisfied. In particular, we modify the Linear Programming (LP) algorithm of [9] by introducing a reduced set of constraints for exact NMF. In contrast to the previous approaches, the proposed algorithm does not require the knowledge of factorization rank (extreme rays [3] or topics [7]). Furthermore, motivated by a similar problem arising in the context of metabolic network analysis [13], we consider an entirely different regime where the number of extreme rays or topics can be much larger than the dimension of the data vectors. The performance of the algorithm for different synthetic data sets are provided.

I. INTRODUCTION

Matrix factorization has numerous applications to the real-world problems where the data matrices representing the numerical observations are often huge and hard to analyze. Meanwhile, factorizing them into lower-rank forms is able to reveal the inherent structure and features, which helps in the meaningful interpretation of the data. In a wide range of natural signals, such as pixel intensities, amplitude spectra, and occurrence counts, negative values are usually physically meaningless. In order to deal with this non-negative constraint, Non-Negative Matrix Factorization (NMF) was introduced.

NMF was first proposed in [1] and used by Lee and Seung for parts-based data representation [2]. It is well known that NMF may not be unique. In this context a sufficient condition on the uniqueness of NMF was pointed out in [3]. A geometrical interpretation [4], of this condition amounts to the fact that the extreme rays (topics) generating the cone (in the non-negative orthant) are contained in the data. Thus, for NMF, one only needs to identify these extreme rays. Additionally, it was demonstrated in [5] that under such a separability assumption, one can use Linear Programming (LP) to isolate the extreme rays from the non-extreme rays. In this paper we will focus only on such cases.

Bittorf et al. [9] presented a LP-based NMF algorithm named *Hottopixx*. Kumar et al. [4], instead, presented a fast conical hull algorithm to deal with the large-scale NMF based on its polyhedral structure. It was shown to perform much faster than *Hottopixx*. However, both of these algorithms require the factorization rank, i.e. the number of extreme rays as a necessary input. Some applications will grant this as prior knowledge but from the view of robustness, the preference would go to a robust NMF algorithm. Gillis and Luce [10] reformulated the algorithm in [9] to detect the number of extreme rays automatically. Nevertheless, the limitation still exists with the fact that the number of constraints in the LP is enormous in face of the large-scale data.

Alternatively, motivated by **Theorem 5.4** in [5], we propose a simpler modification of the constraints to alleviate these issues. In particular we reduce the number of constraints and do not require that the number of extreme rays to be known. This allows us to use a proximal point-based algorithm [12] to solve the LP problem efficiently for data sets with large size.

In addition we also consider an entirely different regime from the NMF applications in the literature so far, where the data lies high dimension with a much smaller factorization rank (the number of extreme rays) in comparison. Specifically, we look at the case when the number of extreme rays is much larger than the dimensionality of the space. This is caused by the computational issues, which arises in Double Description (DD) method [6] for the analysis of metabolic networks to find the Elementary Flux Modes (EFMs) as the set of extreme rays of the polyhedral cone [13]. A NMF like problem arises as an intermediate step in DD method. In this context, we believe that the computational advances in NMF can help with addressing the computational issues in the DD method [6].

The organization of the paper is as follows. Section II provides a brief review of NMF from the geometric perspective as well as the *Hottopixx* algorithm. Section III explains the proposed proximal point algorithm with the reformulated LP constraints. The experiments results are presented in Section IV and the paper concludes in Section V.

Notation: The matrices will be denoted by boldface capital letters and vectors by boldface small letters. In addition we use the MATLAB notation of `diag` to transform vectors to diagonal matrices and to extract the diagonal from the matrix in the argument. Also we use the MATLAB “,” and the “;” operators for matrix concatenations.

II. NON-NEGATIVE MATRIX FACTORIZATION

For the non-noisy case, given a data matrix $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n] \in \mathbb{R}_+^{m \times n}$. Therefore, NMF aims to find two nonnegative matrices $\mathbf{F} \in \mathbb{R}_+^{m \times r}$ and $\mathbf{W} \in \mathbb{R}_+^{r \times n}$ such that $\mathbf{X} = \mathbf{F}\mathbf{W}$. For an approximate NMF, instead, the aim is to solve the following optimization problem.

$$\min_{\mathbf{F}, \mathbf{W} \geq 0} \|\mathbf{X} - \mathbf{F}\mathbf{W}\|_2^2 \quad (1)$$

A. Geometry of the NMF Problem

The factorization $\mathbf{X} = \mathbf{F}\mathbf{W}$ implies that all the *columns* of \mathbf{X} can be represented as non-negative combination of the columns $\{\mathbf{f}_i\}_{i=1}^r$ of the matrix \mathbf{F} . The algebraic characterization can be described as below.

Definition 1. The simplicial cone generated by columns $\{\mathbf{f}_i\}_{i=1}^r$ is given by,

$$\Gamma = \Gamma_F = \{\mathbf{x} : \mathbf{x} = \sum_i \alpha_i \mathbf{f}_i, \alpha_i \geq 0\} \quad (2)$$

The factorization $\mathbf{X} = \mathbf{F}\mathbf{W}$ refers geometrically to that the $\mathbf{x}_i, i = 1, 2, \dots, n$ all lie in or on the surface of the simplicial cone generated by the $\{\mathbf{f}_i\}_{i=1}^r$, as depicted in Fig. 1.

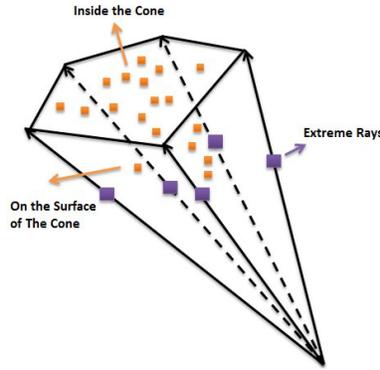


Fig. 1. Geometry of the NMF Problem. Separability implies that data is contained in a cone generated by a subset of r extreme rays (indicated by purple squares).

With this viewpoint in mind, we define three assumptions as follows.

- **Assumption 1:** *Extreme rays* by definition are simplicial: No extreme ray is in the convex combination of the other extreme rays. This is also shown to be necessary and sufficient for exact recovery in topic modeling [8].
- **Assumption 2:** The dataset consisting of all columns of \mathbf{X} , reside in or on the surface of a cone generated by these extreme rays of \mathbf{X} [3].
- **Assumption 3:** Assuming that the columns of \mathbf{X} are normalized to unity there are no duplicate columns in \mathbf{X} .

The above three assumptions will be collectively referred to as *separability assumption* in the following: The entire dataset, i.e. all columns of \mathbf{X} , reside in or on a surface of a cone generated by a small subset of r columns of \mathbf{X} , the vectors in this subset being simplicial and there are no duplicate columns in \mathbf{X} after column normalization.

In algebraic terms, $\mathbf{X} = \mathbf{F}\mathbf{W} = \mathbf{X}_I\mathbf{W}$ for some subset $I \subseteq \{1, 2, \dots, r\}$ of columns (**extreme rays**) of \mathbf{X} and where \mathbf{X}_I denotes the matrix built with columns of \mathbf{X} indexed by I . This means that the r vectors of \mathbf{F} are hidden

among the columns of \mathbf{X} (I is unknown) [5]. Equivalently, it implies that the corresponding subset of r rows of \mathbf{W} constitutes the $r \times n$ weight matrix. Therefore, the computational challenge is to identify the extreme rays efficiently. In this context, we first outline the LP-based *Hottopixx* Algorithm from [9].

B. *Hottopixx*

Bittorf et al. [9] proposed an algorithm of NMF under separability assumption based on the following LP problem:

$$\min_{\mathbf{C} \in \Phi_1(\mathbf{X})} \mathbf{p}^T \text{diag}(\mathbf{C}) \quad (3)$$

and $\mathbf{p} \in \mathbb{R}^{n \times 1}$ is a random vector with distinct positive entries and $\mathbf{C} \in \mathbb{R}_+^{n \times n}$ is referred to as a factorization localizing matrix [9], which belongs to the following polyhedral set.

$$\begin{aligned} \Phi_1(\mathbf{X}) = \{ \mathbf{C} : \mathbf{X}\mathbf{C} = \mathbf{X}, \text{Trace}(\mathbf{C}) = r, \mathbf{C}(i, i) \leq 1 \text{ for all } i \\ \mathbf{C}(i, j) \leq \mathbf{C}(i, i) \text{ for all } i, j, \mathbf{C} \geq 0 \} \end{aligned} \quad (4)$$

For a large scale set-up they proposed an incremental gradient descent algorithm to solve the LP.

III. ROBUST NMF USING *Proximal Point* ALGORITHM

As explained before, two of the prominent shortcomings of existing algorithms for the NMF problem are - (i) Dependence on knowledge of the number of extreme rays r and, (ii) Dealing with a large data set resulting in an enormous number of constraints. An approach in this direction was taken in [10]. However, the number of constraints in their reformulation is still immense for large data. In this paper we present a reformulation which drastically reduces the set of constraints.

A. LP Reformulation

Assuming that the columns of \mathbf{X} are normalized to have an unit \mathbf{L}_1 norm, our LP reformulation for NMF is given as

$$\min_{\mathbf{C} \in \Phi_2(\mathbf{X})} \mathbf{p}^T \text{diag}(\mathbf{C}) \quad (5)$$

where,

$$\Phi_2(\mathbf{X}) = \{ \mathbf{C} : \mathbf{X}\mathbf{C} = \mathbf{X}, \mathbf{C}^T \mathbf{1} = \mathbf{1}, \mathbf{C} \geq 0 \} \quad (6)$$

where $\mathbf{1} \in \mathbb{R}_+^{n \times 1}$ is the vector of all 1-s and $\mathbf{p} \in \mathbb{R}_+^{n \times 1}$ is the same as the vector in (3).

Proposition 1. *Suppose \mathbf{X} admits a separable factorization $\mathbf{F}\mathbf{W}$, compute the minimization of (5) and let $I = \{i : \mathbf{C}_{ii} = 1\}$, then $\mathbf{F} = \mathbf{X}_I$.*

In order to prove the above proposition, we consider the Lagrangian of the optimization problem in (5), which is,

$$\begin{aligned} L(\mathbf{C}, \mathbf{R}, \boldsymbol{\lambda}) = \min_{\mathbf{C}} \mathbf{p}^T \text{diag}(\mathbf{C}) + \text{Tr}\{\mathbf{R}^T(\mathbf{X}\mathbf{C} - \mathbf{X})\} \\ + \boldsymbol{\lambda}^T(\mathbf{C}^T \mathbf{1} - \mathbf{1}) + \text{Tr}\{\mathbf{M}^T \mathbf{C}\} \end{aligned} \quad (7)$$

where \mathbf{R} , $\boldsymbol{\lambda}$ and \mathbf{M} are the Lagrange multipliers. Then the dual form of (5) is

$$\begin{aligned} \max_{\mathbf{R}, \boldsymbol{\lambda}, \mathbf{M}} \quad & - \text{Tr}\{\mathbf{R}^T \mathbf{X}\} - \boldsymbol{\lambda}^T \mathbf{1} \\ \text{s.t.} \quad & \mathbf{P} + \mathbf{X}^T \mathbf{R} + \mathbf{1} \boldsymbol{\lambda}^T + \mathbf{M} = 0, \mathbf{M} \geq 0 \end{aligned} \quad (8)$$

The proof of the proposition follows from **Lemma 1** and **Lemma 2** below.

Lemma 1. *If $\ell \notin I$, $\mathbf{C}_{\ell\ell} = 0$ for all $\mathbf{C} \in \Phi_2(\mathbf{X})$.*

Proof: For $\ell \notin I$, consider the LP problem

$$\min_{\mathbf{C} \in \Phi_2(\mathbf{X})} -\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \quad (9)$$

where $\mathbf{e}_\ell \in \mathbb{R}_+^{n \times 1}$ denotes the vector with ℓ th entry 1 and the rest 0. Assign $\mathbf{P} = -\text{diag}(\mathbf{e}_\ell)$ and using the constraint $\mathbf{C} \geq 0$, we can claim that $-\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \leq 0$. Under the separability assumption, there exists a collection of vectors $\{\boldsymbol{\rho}_i\} \in \mathbb{R}^{n \times 1}$ such that,

$$\begin{aligned} \boldsymbol{\rho}_i^T \mathbf{x}_i &= u_i \\ \boldsymbol{\rho}_i^T \mathbf{x}_j &\leq -v_i, \quad \text{for } j \neq i \end{aligned} \quad (10)$$

for $u_i = 0$ and some $v_i \geq 0$. A feasible solution to (8) is

$$\begin{aligned} \mathbf{P} &= -\text{diag}(\mathbf{e}_\ell), \quad \boldsymbol{\lambda} = \mathbf{0} \in \mathbb{R}^{n \times 1} \\ \mathbf{R} &= [0, \dots, \boldsymbol{\rho}_i, \dots, 0], \quad \text{for some } i \in I \\ \mathbf{M} &= \mathbf{M}_1 + \mathbf{M}_2 : \mathbf{M}_1 = \text{diag}(\mathbf{e}_\ell), \quad \mathbf{M}_2 = -\mathbf{X}^T \mathbf{R} = \mathbf{0} \end{aligned} \quad (11)$$

With such selection, the dual cost function is equal to 0. From *weak duality* [11] it follows that $-\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \geq 0$. Combined with $-\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \leq 0$, it implies $\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) = 0$ and $\mathbf{C}_{\ell\ell} = 0$. ■

Lemma 2. *If $\ell \in I$, $\mathbf{C}_{\ell\ell} = 1$ for all $\mathbf{C} \in \Phi_2(\mathbf{X})$.*

Proof: For $\ell \in I$, Consider the LP problem

$$\min_{\mathbf{C} \in \Phi_2(\mathbf{X})} \mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \quad (12)$$

Note that the constraint $\mathbf{C}^T \mathbf{1} = \mathbf{1}$ implies that $\mathbf{C} \leq 1$ therefore $\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \leq 1$. For the dual program a feasible solution can be found as

$$\begin{aligned} \mathbf{P} &= \text{diag}(\mathbf{e}_\ell), \quad \boldsymbol{\lambda}^T = [0, 0, \dots, -1, \dots, 0] \in \mathbb{R}^{1 \times n}, \quad \text{where } \ell\text{th entry is } -1 \\ \mathbf{R} &= \mathbf{0}, \quad \mathbf{M} = -\mathbf{1}\boldsymbol{\lambda}^T - \mathbf{P} \end{aligned} \quad (13)$$

for which the dual cost function (8) is equal to 1. Again using the *weak duality* [11], it implies that $\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \geq 1$ and from $\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) \leq 1$, we have $\mathbf{e}_\ell^T \text{diag}(\mathbf{C}) = 1$ and $\mathbf{C}_{\ell\ell} = 1$. ■

Proof of Proposition 1. Let \mathbf{C}_0 denote the factorization localizing matrix which identifies the factorization with lowest cost $\mathbf{p}^T \text{diag}(\mathbf{C})$ and has either ones or zeros on the diagonal. Then \mathbf{C}_0 is the unique optimal solution of (5). To see this, let I denote the set of simplicial columns of minimum cost. Since each column of \mathbf{X} can only belong to I or not, once \mathbf{X} is given, \mathbf{C} is determined then the lowest cost $\mathbf{p}^T \text{diag}(\mathbf{C})$ is determined, which is unique.

Remark. *Note that \mathbf{W} can be readily obtained from \mathbf{C} as $\mathbf{W} = \mathbf{C}(I, :)$ (in MATLAB notation).*

B. Proximal Point Algorithm

Based on the above reformulation, a proximal point-based (*Proximal Point*) algorithm is employed in this paper. A necessary pre-processing step is the column normalization, which makes sure that the sum of each column of \mathbf{X} is equal to one.

After the normalization, we can rewrite the LP in (5) as

$$\min_{\mathbf{C} \in \Phi_3(\mathbf{X})} \mathbf{p}^T \text{diag}(\mathbf{C}) \quad (14)$$

where,

$$\Phi_3(\mathbf{X}) = \{\mathbf{C} : \mathbf{A}\mathbf{C} = \mathbf{A}, \mathbf{C} \geq 0\} \quad (15)$$

where $\mathbf{A} = [\mathbf{X}; \mathbf{1}^T]$

We solve this LP using the proximal point algorithm [12] in **Algorithm 1** and in our implementation, $\{t_k\}$ are set to a large constant 100. The discussion on the convergence of the algorithm can be found in [12].

IV. EXPERIMENTS RESULTS

All of the experiments were run on an identical configuration: a dual Xeon W3505 (2.53GHz) machine with 6GB RAM. *Proximal Point Algorithm* is examined in MATLAB with the version of 2013a.

Algorithm 1 Robust NMF by *Proximal Point* Algorithm

Input: A column normalized matrix $\mathbf{X} \in \mathbb{R}_+^{m \times n}$, stopping threshold ϵ .

Output: A matrix $\mathbf{F} \in \mathbb{R}_+^{m \times r}$ and $\mathbf{W} \in \mathbb{R}_+^{r \times n}$, and $\mathbf{X} = \mathbf{F}\mathbf{W}$.

1: Initialize $\mathbf{Q}^0 = 0$ and $\mathbf{C}^0 = 0$, randomly generate $\mathbf{p} \in \mathbb{R}_+^{m \times 1}$.

2: Update \mathbf{C}^{k+1} :

$$\begin{aligned} \mathbf{C}^{k+1} &= \underset{\mathbf{C}}{\operatorname{argmin}} \{ \mathbf{p}^T \operatorname{diag}(\mathbf{C}) + \frac{1}{t^k} \|\mathbf{Q}^k + t^k(\mathbf{A}\mathbf{C} - \mathbf{A})\|_2^2 \} \\ &= \frac{1}{2t^k} (\mathbf{A}^T \mathbf{A})^{-1} \left(2t^k \mathbf{A}^T \mathbf{A} - \operatorname{diag}(\mathbf{p}) - 2\mathbf{A}^T \mathbf{Q}^k \right) \end{aligned}$$

3: Project \mathbf{C}^{k+1} to the constraint $\mathbf{C} \geq 0$ using $\mathbf{C}^{k+1} = \operatorname{pos}(\mathbf{C}^{k+1})$, where $\operatorname{pos}(\cdot)$ keeps the positive elements and switch the negative elements to 0.

4: Update \mathbf{Q}^{k+1} : $\mathbf{Q}^{k+1} = \mathbf{Q}^k + t^k(\mathbf{A}\mathbf{C}^{k+1} - \mathbf{A})$.

5: Stop the iterations if $\|\mathbf{C}^{k+1} - \mathbf{C}^k\|_2 \leq \epsilon$.

6: Let $I = \{i : C_{ii} = 1\}$ and set $\mathbf{F} = \mathbf{X}_I$ as well as obtain $\mathbf{W} = \mathbf{C}(I, :)$.

TABLE I
EXPERIMENTS ON DIFFERENT DATASET REGARDING TO C1-C3

Data Set	# of Extreme Rays	Accuracy	ϵ
100 × 75(C1)	25	25/25	10 ⁻⁵
500 × 375(C1)	25	23/25	10 ⁻⁴
1200 × 600(C1)	300	300/300	10 ⁻⁴
25 × 100(C2)	15	14/15	10 ⁻⁵
125 × 500(C2)	75	74/75	10 ⁻⁴
425 × 1200(C2)	225	223/225	10 ⁻⁴
25 × 100 (C3)	45	45/45	10 ⁻⁵
125 × 500 (C3)	150	150/150	10 ⁻⁴
425 × 1200(C3)	625	625/625	10 ⁻⁴

A. Random Data Generation

To generate our instances, r independent extreme rays are firstly created in $\mathbb{R}_+^{m \times 1}$, with the element value between $[0, 100]$. The remaining columns are then generated to be the random non-negative combinations of the r' extreme rays, where $r' \in [2, r]$ is randomly selected for each of the $n - r$ points. The column normalization is carried out sequentially. Three regimes of NMF problems are analyzed here:

- (C1). $m \geq n, m \geq r$, which is motivated from the data structure for face recognition [14]
- (C2). $r \leq m \leq n$, which is the scenario for topic modeling problem [7]
- (C3). $m \leq r \leq n$, which can be applied to metabolic network data [13].

Furthermore, since the algorithm is free from the order of the columns, the r extreme rays are allocated at the beginning of each data set.

Different size of data sets are generated to check the effectiveness of *Proximal Point* Algorithm, from small to large-scale. In Tab. I, the last column indicates the highest level for iteration stopping criterion ϵ to achieve the listed accuracy. From the experiments, it is exhibited that our algorithm can deal with three regimes of the data with different sizes. Moreover, the identification accuracy is satisfying.

B. Application to Image Processing

In this section, we apply the *Proximal Point* algorithm to one face image processing data set, namely, CBCL Dataset [14]. Basically, the CBCL face dataset is made of 2429 gray-level face images with 19×19 pixels. We randomly choose 20 images from the dataset with vectorization to be the generators, which means the number of extreme rays in this case is $r = 20$. Through the random non-negative combination of the extreme rays, a 361×500

facial data matrix is created. Applying *Proximal Point* algorithm to this dataset, the results of the extreme rays identification are shown in Fig. 2, which represents the initial 20 images as generators.

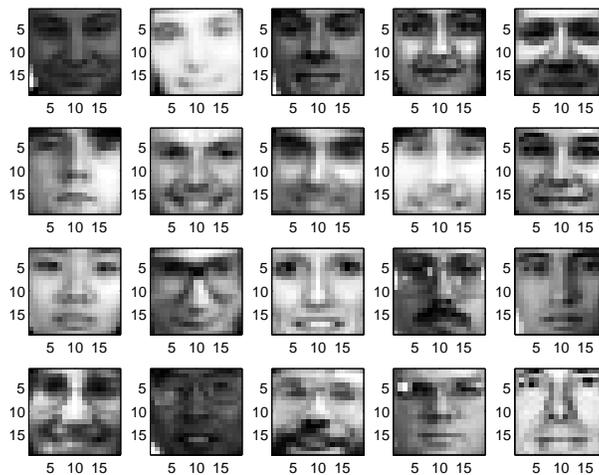


Fig. 2. The facial images identified as extreme rays for the 361×500 data set with 20 extreme rays. The stopping criterion was selected as 10^{-5} .

V. ACKNOWLEDGEMENTS

The second author would like to acknowledge several useful discussions with Prof. Prakash Ishwar at Boston University.

REFERENCES

- [1] P. Paatero and U. Tapper, "Positive Matrix Factorization: A non-negative factor model with optimal utilization of error estimates of data values," *Environmetrics*, vol. 5, pp. 111-126, 1994.
- [2] D. Lee and H. Seung, "Algorithms for Non-Negative Matrix Factorization," in *NIPS*, 2001, pp. 556-562.
- [3] D. Donoho and V. Stodden, "When does non-negative matrix factorization give a correct decomposition into parts, in *NIPS*, 2004, MIT Press.
- [4] A. Kumar, V. Sindhwani, and P. Kambadur, "Fast conical hull algorithms for near-separable non-negative matrix factorization," arXiv:1210.1190, 2012.
- [5] S. Arora, R. Ge, R. Kannan, and A. Moitra, "Computing a nonnegative matrix factorization – provably," in *STOC*, 2012, pp.145-162.
- [6] K. Fukuda and A. Prodon. "Double description method revisited." in *Combinatorics and Computer Science*, vol. 1120, pp. 91-111, 1996.
- [7] W. Ding, M. H. Rohban, P. Ishwar, V. Saligrama, "A New Geometric approach to latent topic modeling and discovery," arXiv:1301.0858 [stat.ML], 2013.
- [8] W. Ding, P. Ishwar, M. H. Rohban, V. Saligrama, "Necessary and Sufficient Conditions for Novel Word Detection in Separable Topic Models," arXiv:1310.7994 [cs.LG], 2013.
- [9] B. Recht, C. Re, J. Tropp, and V. Bittorf, "Factoring nonnegative matrices with linear programs," in *NIPS*, 2012, pp. 1223-1231.
- [10] N. Gillis and R. Luce, "Robust near-Separable nonnegative matrix factorization using linear optimization," arXiv:1302.4385, 2013.
- [11] D.G. Luenberger, *Optimization by vector space methods*. New York: Wiley, 1969.
- [12] J. Eckstein, "Nonlinear proximal point algorithms using Bregman functions with applications to convex programming," *Math. Oper. Res.*, vol. 18, pp. 203-226, 1993.
- [13] M. Terzer and J. Stelling, "Accelerating the Computation of Elementary Modes Using Pattern Trees," in *WABI*, 2006, pp. 333-343.
- [14] <http://cbcl.mit.edu/software-datasets/FaceData2.html>
- [15] D. P. Bertsekas, "A new class of incremental gradient methods for least squares problems," *SIAM J. Optim.*, vol. 7, no. 4, pp. 913-926, 1997.